

# 1 Super-resolution fight club: A broad assessment of 2D & 3D 2 single-molecule localization microscopy software

3 *Daniel Sage*<sup>\*+1</sup>, *Thanh-An Pham*<sup>+1</sup>, *Hazen Babcock*<sup>2</sup>, *Tomas Lukes*<sup>3</sup>, *Thomas Pengo*<sup>4</sup>, *Ramraj Velmurugan*<sup>5</sup>, *Alex*  
4 *Herbert*<sup>6</sup>, *Anurag Agrawal*<sup>7</sup>, *Silvia Colabrese*<sup>1,8</sup>, *Ann Wheeler*<sup>9</sup>, *Anna Archetti*<sup>10</sup>, *Bernd Rieger*<sup>11</sup>, *Raimund Ober*<sup>5</sup>,  
5 *Guy M. Hagen*<sup>12</sup>, *Jean-Baptiste Sibarita*<sup>13</sup>, *Jonas Ries*<sup>14</sup>, *Ricardo Henriques*<sup>15</sup>, *Michael Unser*<sup>1</sup>, *Seamus Holden*<sup>\*+16</sup>

6 \*Corresponding authors: [daniel.sage@epfl.ch](mailto:daniel.sage@epfl.ch), [seamus.holden@ncl.ac.uk](mailto:seamus.holden@ncl.ac.uk).

7 +Equal contribution

8 1: Biomedical Imaging Group, School of Engineering, Ecole Polytechnique Fédérale de Lausanne  
9 (EPFL), Switzerland

10 2: Harvard Center for Advanced Imaging, Harvard University, Cambridge, Massachusetts, USA

11 3: Laboratory of Nanoscale Biology & Laboratoire d'Optique Biomédicale, STI - IBI, EPFL, Lausanne,  
12 Switzerland

13 4: University of Minnesota Informatics Institute, University of Minnesota Twin Cities, USA

14 5: Electrical Engineering, University of Texas Dallas, Richardson, Texas, USA

15 6: MRC Genome Damage and Stability Centre, School of Life Sciences, University of Sussex, Brighton,  
16 UK

17 7 : Double Helix LLC, Boulder, Colorado, USA

18 8 : Istituto Italiano di Tecnologia, Genova, Italy

19 9: Edinburgh Super-Resolution Imaging Consortium, University of Edinburgh, UK

20 10 : Laboratory of Experimental Biophysics, École Polytechnique Fédérale de Lausanne (EPFL),  
21 Lausanne, Switzerland

22 11: Department of Imaging Physics, Faculty of Applied Sciences, Delft University of Technology, The  
23 Netherlands

24 12: UCCS center for the Biofrontiers Institute, University of Colorado at Colorado Springs, Colorado,  
25 USA

26 13: Institut Interdisciplinaire de Neurosciences, University of Bordeaux, France

27 14: European Molecular Biology Laboratory (EMBL), Cell Biology and Biophysics Unit, Heidelberg,  
28 Germany

29 15: Quantitative Imaging and Nanobiophysics Group, MRC Laboratory for Molecular Cell Biology,  
30 University College London, UK

31 16: Centre for Bacterial Cell Biology, Institute for Cell and Molecular Biosciences, Newcastle  
32 University, UK

33

34

## 35 **ABSTRACT**

36 With the widespread uptake of 2D and 3D single molecule localization microscopy, a large set of  
37 different data analysis packages have been developed to generate super-resolution images. To guide  
38 researchers on the optimal analytical software for their experiments, we have designed, in a large  
39 community effort, a competition to extensively characterise and rank these options. We generated  
40 realistic simulated datasets for popular imaging modalities – 2D, astigmatic 3D, biplane 3D, and double  
41 helix 3D – and evaluated 36 participant packages against these data. This provides the first broad  
42 assessment of 3D single molecule localization microscopy software, provides a holistic view of how  
43 the latest 2D and 3D single molecule localization software perform in realistic conditions, and  
44 ultimately provides insight into the current limits of the field.

## 45 INTRODUCTION

46 Image processing software is central to single molecule localization microscopy (SMLM), which  
47 delivers an order of magnitude resolution improvement on diffraction limited conventional  
48 fluorescence microscopy, from 250 nm to approximately 20 nm resolution, by temporal separation of  
49 fluorophores within a sample<sup>1-3</sup>. Efficient and automated image processing is essential to extract the  
50 super-resolved positions of individual molecules from thousands of raw microscope images,  
51 containing millions of blinking fluorescent spots.

52 Improvements in SMLM image processing algorithms have been crucial in maximizing spatial  
53 resolution and in reducing the imaging time of SMLM for compatibly with live cell imaging<sup>4-6</sup>. If SMLM  
54 is to achieve a resolving power approaching that of electron microscopy, the analysis software  
55 employed needs to be robust, accurate, and performing at current algorithmic limits. This can only be  
56 achieved through rigorous quantification of SMLM software performance.

57 The first localization microscopy software challenge was carried out in 2013, to enable robust  
58 benchmarking of 2D localization microscopy software packages<sup>7</sup>. But biology is not just a 2D problem,  
59 and a key focus of localization microscopy is the imaging of 3D imaging of nanoscale cellular  
60 processes<sup>8,9</sup>. 3D localization microscopy is a more difficult image processing problem than 2D SMLM.  
61 In addition to finding the center of diffraction limited spots to super-resolve lateral position, 3D SMLM  
62 algorithms must also extract axial information from the image, usually by measuring small changes in  
63 the shape of a fluorophore's PSF<sup>10</sup>.

64 There are roughly three common approaches for 3D SMLM. First, point spread function engineering,  
65 where the axial asymmetry of the microscope point spread function (PSF) is increased by introducing  
66 intentional aberrations in the system, ranging from simple astigmatism<sup>10</sup> to more complex PSF  
67 manipulation such as the double helix PSF method<sup>11</sup>. Second, biplane or multiplane imaging, where  
68 axial position is measured based on simultaneous measurement of PSF shape at two or more focal  
69 planes<sup>12</sup>. Third, dual objective based interferometry, where Z-position is calculated from single photon  
70 interference between opposing objectives<sup>13</sup>. Multiplane and PSF engineering methods typically obtain  
71 axial resolutions on the order of 50 nm<sup>10,11</sup>. Interferometry achieves the best axial resolution, 10-20  
72 nm<sup>13</sup>, but is not yet widely adopted.

73 Despite the widespread use of 3D localization microscopy, and challenging nature of 3D SMLM image  
74 processing, the performance of software for 3D single molecule localization microscopy has previously  
75 only been assessed for 2 or 3 software packages at a time, and without standard test data or metrics<sup>14-</sup>  
76 <sup>17</sup>. In the absence of common reference datasets and reliable assessment procedure of 3D software  
77 performance, it is not possible to objectively assess how different software affects final image quality,  
78 or which algorithmic approaches are most successful. Crucially, end-users cannot determine which 3D  
79 SMLM software package and imaging modality is optimal for their application.

80 We therefore ran the first 3D localization microscopy software challenge, to assess the performance  
81 of 3D SMLM software. We generated synthetic datasets for three popular 3D SMLM modalities:  
82 astigmatic imaging, biplane imaging and double helix point spread function microscopy. We also ran  
83 a second 2D localization microscopy software challenge, to reassess the 2D SMLM software state-of-  
84 the art on new, tougher, more realistic datasets.

85 Our simulations incorporate experimentally acquired point spread functions for maximal authenticity,  
86 used signal and noise levels based closely on common experimental conditions, and incorporated a  
87 realistic 4-state model of fluorophore photophysics<sup>18</sup>. Our synthetic data was designed to mimic two  
88 common classes of cellular structure: narrow line-like microtubules (MT) and larger tubes similar to  
89 the endoplasmic reticulum (ER) or mitochondria. Our simulations also included conditions with low  
90 density (LD) of active fluorophores, used experimentally to obtain maximal resolution, and with high  
91 density (HD) of active fluorophores, used experimentally for fast or live cell imaging.

## 92 RESULTS

### 93 Competition design

94 We established a large committee from within the SMLM research community, including  
95 experimentalists and software developers, to define the scope of the challenge, ensure realism of the  
96 datasets and define analysis metrics. We further opened this discussion to the whole community,  
97 through an open forum, discussing best practices for the implementation of this contest<sup>19</sup>.

98 Thirty-six software packages have been entered in the competition thus far (Table S1). Excitingly,  
99 participation in the competition actually led at least 8 teams to their software to support additional  
100 3D SMLM modalities, showing how competition fosters microscopy software development.

101 In 2016, we ran a first round of the 3D SMLM competition with explicit submission deadlines, with 30  
102 competitor teams, culmination in a special session at the 6th annual Single Molecule Localization  
103 Microscopy Symposium (SMLMS 2016). Since then, the challenge has been opened to continuously  
104 accept new entries. We have had 12 new registrations of which 5 have submitted localizations,  
105 including a multiple best-in-class performer (SMAP-2018<sup>20</sup>, an updated version of previously entered  
106 software) demonstrating the utility of the competition as an evolving measure of the state of the field.

### 107 Realistic 3D simulations

108 Testing super-resolution software on experimental data lacks the ground truth information required  
109 for rigorous quantification of software performance. Therefore, realistic simulated 3D SMLM datasets  
110 are required. After comparison of simulated microscope PSFs with multiple experimental PSFs from  
111 SMLM microscopes around the world, we observed that a critical challenge to realistic 3D SMLM  
112 simulations was to accurately model the experimental microscope PSF for each 3D modality. Even  
113 experimental 2D PSFs showed significant aberrations away from the focal plane (Fig S9).

114 3D SMLM inherently involves addition of aberrations to the microscope PSF to encode the Z-position  
115 of the molecule. For the PSF models included in the competition: 2D, astigmatic (AS), double helix  
116 (DH), and biplane (BP), we observed that the PSFs showed complex aberrations not well described by  
117 simple analytical models (Fig S9). We thus combined experimental 3D PSFs with simulated ground  
118 truth by performing simulations using PSFs directly derived from experimental calibration data (Fig 1,  
119 *Methods*). The experimental PSFs used to generate the simulated data are available online (*Methods*)  
120 and are representative of 3D SMLM PSFs obtained on typical microscopes.

121 For the 3D competition, we simulated synthetic 25 nm diameter microtubules (Fig 1). For the 2D  
122 competition, in addition to synthetic microtubules (MT), we simulated larger diameter 150 nm  
123 cylinders, designed to approximate larger cellular structures such as mitochondria and the  
124 endoplasmic reticulum (ER) (Fig 1). We incorporated a 4-state model of fluorophore photophysics,  
125 including a transient dark state (dye “blinking”) and a bleaching pathway (Fig S1C).

126 As performance at different density of active emitters is a key challenge for SMLM software, we  
127 generated 3D competition datasets at both sparse emitter density (0.2 mol. [molecule]  $\mu\text{m}^{-2}$ ) and high  
128 emitter density (2 mol.  $\mu\text{m}^{-2}$ ). We additionally generated a very high density dataset (5 mol.  $\mu\text{m}^{-2}$ ) for  
129 the 2D competition.

130 We generated data at three different signal-to-noise ratio (SNR) levels, based on real signal to noise  
131 levels encountered under common SMLM experimental scenarios: fixed cells antibody labelled with  
132 organic dye<sup>10</sup>, fluorescent protein labelling<sup>1</sup>, and live cell affinity dye labelling<sup>21,22</sup>.

133 Together, these simulations closely resemble experimental 3D and 2D data under a range of  
134 challenging conditions of SNR, spot density, axial thickness and test structure summarized in Table S2.  
135 In addition, we provide also a z-stack of extremely bright beads for software calibration. The  
136 competition datasets are available online (*Methods*).

## 137 **Quantitative performance metrics for comparison of 3D software**

138 We assessed software performance by 24 quality metrics (*Supplementary Note 2*), in four categories:  
139 1) single molecule localization error, 2) ability to successfully detect molecules, 3) image-based  
140 resolution metrics and 4) image-based signal to noise ratio (*Methods*). We also recorded software run  
141 time. The complete set of summary statistics, axially resolved performance and super-resolved images  
142 is available for each competition software on the competition website. We generated an online  
143 leaderboard<sup>23</sup>, allowing easy ranking of each software by each metric. Software results can be  
144 accessed interactively through a visualization interface that allows side-by-side comparison of results  
145 for multiple software packages (Fig S10).

146 In order to rank overall software performance, we performed a principal component analysis of core  
147 metrics to identify key variables (Fig S14). A correlation matrix identified four major blocks of metrics  
148 showing strong codependency (Fig S14B) corresponding closely to the manually identified categories  
149 above. We chose to focus further analysis primarily on the metrics directly derived from single  
150 molecule localizations, rather than image derived metrics, which we reasoned would be sensitive to  
151 additional factors such as image rendering method. We thus chose representative metrics from the  
152 first two blocks:

153 *1. Single molecule localization error.* The foremost consideration for localization software is how  
154 accurately it finds the position of labelled molecules. This was quantified as the root mean squared  
155 localization error (RMSE) between the measured molecule position and the ground truth, in both the  
156 lateral (XY) and axial (Z) dimensions.

157 *2. Ability to successfully detect fluorescent molecules.* In addition to localization precision, SMLM  
158 image resolution also depends critically on number of localized molecules<sup>24</sup>, so it is crucial for SMLM  
159 software to accurately detect a large fraction of molecules in a dataset, and minimize false  
160 localizations. For every frame, we identified the localizations that are close enough to a ground-truth  
161 position as true-positives (TP), the spurious localizations as false-positives (FP) and the undetected  
162 molecules as false-negatives (FN). We then computed the *Jaccard index* (JAC, %), which measures the  
163 fraction of correctly detected molecules in a dataset:

$$164 \quad JAC = 100 \frac{TP}{TP + FP + FN}$$

165 The average JAC, lateral RMSE and axial RMSE measured the performance of a software. A very good  
166 RMSE should always read in context of the Jaccard index to check if good RMSE is not obtained only  
167 for the brightest molecules.

168 For ranking purpose, we developed a single summary statistic for overall evaluation of software  
169 performance, which we term the *efficiency* (*E*), encapsulating both the software's ability to find  
170 molecules, measured by the Jaccard index, and the software's ability to precisely localize molecules.

$$171 \quad E = 100 - \sqrt{(100 - JAC)^2 + \alpha^2 \cdot RMSE^2}$$

172 The trade-off between these two metrics is controlled by a parameter  $\alpha$ . In a retrospective analysis,  
173 we chose  $\alpha = 1 \text{ nm}^{-1}$  for the lateral efficiency  $E_{\text{lat}}$ ,  $\alpha = 0.5 \text{ nm}^{-1}$  for the axial efficiency  $E_{\text{ax}}$ , based on the  
174 linear regression slope between the localization errors and Jaccard index (Fig 14A). Using this  
175 definition, an average software performance has an efficiency in the range 25-75, ground-truth has  
176 the maximum efficiency of 100. Overall 3D efficiency was calculated as the average of lateral and axial  
177 efficiencies.

## 178 **Performance of 3D software**

179 Complete rankings for each imaging modality and spot density are presented (Fig 2, S13), together  
180 with summary information on all competition software (Table S1, *Supplementary Note 1*). As these  
181 data are continuously updated on the competition website, this resource provides microscopists with

182 an easy quick reference for the current state of the art, including current best-in-class performers for  
183 each category.

184 After assembling an overall summary of best performers for each competition category, we  
185 investigated the performance of software within each imaging modality.

### 186 *Astigmatic localization microscopy*

187 Astigmatic localization microscopy is probably the most popular imaging 3D SMLM modality, reflected  
188 by the highest number of software submissions in the 3D competition (Fig 2). For astigmatism, we  
189 observed a large spread of software performance, even for the most straightforward high SNR, low  
190 spot density (LD) conditions (Fig 3A-B, Table S5). The best-in-class software (SMAP-2018) has  
191 significantly better localization error and Jaccard index performance than average (lateral RMSE 26  
192 nm best vs 38 nm average, axial RMSE 29 nm best vs 66 nm average, Jaccard index 85 % best vs 74 %  
193 average). Clearly, the quality of the image reconstruction depends strongly on choice of 3D software.

194 To investigate the reasons for software variation, we inspected plots of software performance as a  
195 function of axial position in the low density, high SNR dataset for best-in-class and representative  
196 middle-range software (Fig S6A). We observed that the key cause of the spread in software  
197 performance is variation in software performance away from the focal plane. Near the focal plane,  
198 most software packages perform well. However, the axial and lateral RMSE away from the plane of  
199 focus is significantly higher for the best in class software, and the Jaccard index is also slightly improved  
200 (Fig 6A). This is also visibly apparent in the super-resolved images (Fig 4, top panel). We observed that  
201 best-in-class software had a Z-range (the FWHM range of axially resolved software recall, *Methods*) of  
202 1170 nm, greater than two-thirds of the simulated range. Outside this range, the recall and Jaccard  
203 index dropped sharply, probably due the large increase in PSF size and decrease in effective SNR at  
204 significant defocus (Fig S9).

205 When we examined results for the low SNR, low density dataset (Fig 2B, 3B), we found an expected 2-  
206 fold degradation in best-in-class RMSE (lateral RMSE 39 nm, axial RMSE 60 nm), due to the decrease  
207 in image SNR. However, the best-in-class software (SMolPhot) Jaccard index was effectively constant  
208 between the low and high SNR datasets (86 % vs 85 %), although the Z-range did drop at lower SNR  
209 (930 nm vs 1120 nm). The best astigmatism software packages were thus remarkably good at finding  
210 spots at low SNR, even away from the plane of focus.

211 We analyzed how close software performance was to theoretical limits by calculating the Cramer-Rao  
212 Lower Bound (CRLB) as a function of axial position for each dataset and comparing it to the best-in-  
213 class software results (Fig S7, Fig S8). Close to the focus, best-in-class software was close to CRLB  
214 performance, but significant deviations for the CRLB limit occurred > 200 nm. This could be due to the  
215 difficulty in actually detecting the spots away from focus.

216 When we examined astigmatic software performance for the challenging high spot density datasets  
217 (Fig 2B, 3), performance was reduced. For the high SNR high spot density dataset (best software,  
218 SMolPhot), localization error increased and Jaccard index decreased significantly compared to the low  
219 density condition (lateral RMSE best HD 51 nm vs best LD 27 nm, axial RMSE best HD 66 nm vs best  
220 LD 29 nm, Jaccard index best HD 66 % vs best LD 85 %). Inspection of the super-resolved images (Fig 4)  
221 nevertheless shows acceptable results for the HD dataset, particularly in the lateral dimension. In  
222 many circumstances, the performance reduction at 10x higher spot density should be acceptable for  
223 10x faster, potentially live-cell-compatible, imaging speed. We also observed a large spread of  
224 software performance for the high density datasets, probably because a significant fraction of the  
225 software packages were primarily designed for low density conditions.

226 We observed poor performance for the most challenging low SNR high spot density astigmatism  
227 dataset (Fig 2, 3, S3, best software SMolPhot). Best-in-class localization precision and Jaccard index  
228 decreased significantly (lateral RMSE 76 nm, axial RMSE 101 nm, Jaccard index 58 %). These data

229 suggest that low SNR high density 3D astigmatic localization microscopy entails a significant reduction  
230 in image resolution.

### 231 *Double helix point spread function localization microscopy*

232 We next analyzed the performance of the double helix software (Fig S13). For the software in the high  
233 SNR low spot density condition, double helix software showed more uniform performance than  
234 astigmatism. Best-in-class software (SMAP-2018) showed only a limited improvement compared with  
235 average software (Fig 3B, lateral RMSE, 27 nm best vs 37 nm average; axial RMSE 21 nm best vs 34 nm  
236 average; Jaccard index 77 % best vs 73 % average). In general software localization performance was  
237 close to the CRLB (Fig S7, S8). We observed that performance of the software away from the focal  
238 plane is relatively uniform (Fig 4, S6A), and best-in-class Z-range at high SNR was large at 1180 nm (Fig  
239 S6). Double helix imaging may show less software-to-software variation and large Z-range at low spot  
240 density than astigmatic imaging because the PSF shape and intensity are fairly constant as a function  
241 of Z – compared to astigmatic imaging, where spot size, shape and intensity vary greatly as a function  
242 of Z (Fig S9).

243 Double helix software performance decreased significantly for the low spot density low SNR condition  
244 (best software SMAP-2018), particularly in terms of best-in-class Jaccard index (66 % low SNR vs 77 %  
245 high SNR, Figure 3B, S3, S13A). DH Jaccard index was also significantly worse than astigmatism results  
246 at either high or low SNR (85 % high SNR, 86 % low SNR). This indicates that it was quite hard to  
247 successfully find localizations in the low SNR DH dataset, likely because the large size of the DH PSF  
248 spreads emitted photons over a large area, lowering effective image SNR.

249 Double helix software performed poorly on the high spot density datasets at high SNR (best software  
250 CSpline), especially in terms of the Jaccard index (Fig 3B, S13A, best lateral RMSE 67 nm, best axial  
251 RMSE 69 nm, best Jaccard index 46 %). The poor performance at high spot density is again probably  
252 because the large DH PSF size increases spot density and decreases SNR (Fig S9). DHPSF performance  
253 at high spot density and low SNR was also not reliable (Fig. 3B, S13A, best software SMAP-2018).

### 254 *Biplane localization microscopy*

255 Best-in-class biplane software (SMAP-2018), at low spot density and for both high and low SNR,  
256 delivered the best performance in any modality (high SNR: lateral RMSE 12.3 nm, axial RMSE 21.7 nm,  
257 Jaccard 87 %), despite a slightly decreased image SNR for the biplane simulations (*Methods*). We  
258 observed a significant spread in software performance in terms of lateral RMSE and Jaccard index,  
259 with the best-in-class software significantly outperforming the other competitors (Fig S13B, 2D). At  
260 low spot density, best-in-class biplane software (SMAP-2018) showed good performance as a function  
261 of Z, with high Jaccard index over almost the entire Z-range of the simulations, and with a Z-range of  
262 1200 nm at high SNR (Fig S6A, C, Table S5). The axial RMSE was relatively uniform as a function of Z  
263 and close to the CRLB limit (Fig S7). As axial and lateral RMSE are both averaged over the entire Z-  
264 range, the strong biplane results arise from good performance across a large Z-range (Fig S6).

265 At high spot density and high SNR, best-in-class biplane software (SMAP-2018) showed acceptable  
266 super-resolved performance (Fig 3B, 4, S13B, best lateral RMSE 43 nm, best axial RMSE 49 nm, best  
267 Jaccard index 61 %). Uniquely among the 3D modalities, best-in-class biplane software also gave  
268 acceptable performance at high spot density and low SNR (Fig 3B, 4, S13B, best lateral RMSE 55 nm,  
269 best axial RMSE 72 nm, best Jaccard index 61 %, best software SMAP-2018).

## 270 **Performance of 2D software**

271 Alongside the 3D challenge, we ran a second edition of the 2D localization microscopy software  
272 challenge<sup>7</sup> to assess how the latest 2D software performed on more challenging, more realistic  
273 datasets, and to provide an assessment of how the field had progressed since the last challenge. We  
274 used the new simulation software, including an experimentally derived PSF and a realistic blinking  
275 model, and also simulated a very high spot density condition (5 molecules/ $\mu\text{m}^2$ ). We created a more

276 spatially extended test structure, "pseudo-endoplasmic reticulum" (pseudo-ER), composed of 150 nm  
277 diameter hollow tubes, to avoid artefacts due to 1D simulated structures<sup>25</sup>. We generated two  
278 different imaging conditions with overall similar SNR but different brightness properties; one with low  
279 fluorophore brightness and low autofluorescence (the low SNR condition for the 3D challenge,  
280 designed to simulate fluorescent protein based SMLM, Fig S4) and one with high fluorophore  
281 brightness and high autofluorescence (to simulate affinity-dye-based live cell SMLM, Fig S5). We used  
282 lateral RMSE, Jaccard index and overall lateral efficiency to rank the 2D software (Fig 2, S2, Table S1).

283 For the pseudo-ER dataset, at low density, best-in-class software (ADCG) performed well (Fig. S4, S5),  
284 with a Jaccard index of 90 % and lateral RMSE of 31 nm, substantially better than the class average  
285 (Jaccard index 72 %, lateral RMSE 36 nm). Low density results for the dimmer fluorophore  
286 microtubules dataset were similar to the brighter pseudo-ER dataset (Fig S2, best software SMolPhot).  
287 For the very high density 2D dataset, which had 25x higher spot density than the LD dataset, best-in-  
288 class software (ADCG) showed excellent performance, with Jaccard index of 75% and lateral RMSE of  
289 45.5 nm (Fig S2). Best-in-class performance (ADCG) on the dimmer fluorophore data at high spot  
290 density was also strong (Fig S2, best Jaccard index 70 %, best lateral RMSE 51 nm).

## 291 Algorithms

292 We identified several classes of algorithm participant software (Table S1):

293 1) *Non-iterative* software tends to regroup the pixels in the local neighborhood of the candidates, like  
294 interpolation, center of mass (QuickPALM<sup>26</sup>) or template matching (WTM<sup>27</sup>). These (often older)  
295 algorithms are fast but tend to achieve poor performance (Table S1).

296 2) *Single emitter fitting* software is usually built on a multi-step strategy of detection, spot localization,  
297 and optional spot rejection. The detection step finds bright spots in noisy images on the pixel grid. The  
298 selection of candidates is usually performed by local maximum search after a denoising filter. Others  
299 rely on more complex algorithms like the wavelet transform (*e.g.*, WaveTracer<sup>28</sup>). We did not observe  
300 software ranking to depend significantly on the choice of optimization scheme, least-square, weighted  
301 least-square or maximum-likelihood estimator (Table S1).

302 3) *Multi-emitter fitting* software groups clusters of overlapping spots, and simultaneously fits multiple  
303 model PSFs to the data. Typically, fitted spots are added to the cluster until a stopping condition is  
304 met<sup>4,5</sup>. This leads to improved localization performance at high spot density, at the cost of reduced  
305 speed. This class of software (*e.g.*, 3D-DAOSTORM<sup>14</sup>, CSpline<sup>14</sup>, PeakFit, ThunderSTORM<sup>29</sup>) was  
306 amongst the top performers in each 2D and 3D competition category (Table S1).

307 As expected, single- and multiple-emitter fitting methods both performed well on low density data  
308 (Table S1); apparently at the densities studied, exclusion of occasionally overlapping spots by single-  
309 emitter software is sufficient for strong performance; explicit multi-emitter fitting is not required. For  
310 the 2D challenge, multi-emitter fitting showed a clear advantage over single emitter fitting at high  
311 density (Table S1). Surprisingly however, well-tuned single-emitter fitting algorithms (SMolPhot,  
312 SMAP-2018) outperformed multi-emitter algorithms for the 3D high density conditions.

313 4) *Compressed sensing algorithms*. One subset of these algorithms utilize deconvolution with sparsity  
314 constraints to reconstruct super-resolved images<sup>30-32</sup>. Although deconvolution approaches can give  
315 good results, they are limited by the necessary use of a sub-pixel grid; increased localization precision  
316 requires smaller grid resolution, which must be balanced against increased computational time.  
317 Recent approaches address this issue by localizing the point sources in a grid-less manner using an  
318 alternating descent conditional gradient scheme under some sparsity constraint (ADCG<sup>33</sup>, SMfit,  
319 SOLAR\_STORM, TVSTORM<sup>34</sup>). This software class consistently gave the overall best performance for  
320 2D high-density (ADCG<sup>33</sup> 1<sup>st</sup>, FALCON<sup>32</sup> 2<sup>nd</sup>, SMfit 3<sup>rd</sup>).

321 5) *Other approaches*. Of the alternative algorithmic approaches used (Table S1), the annihilating filter-  
322 based method LEAP<sup>35</sup> gave good performance for biplane imaging (the only condition for which it was  
323 entered).

#### 324 *Post-hoc temporal grouping*

325 Because molecule on-time is stochastically distributed across multiple frames, a common post-  
326 processing approach to improve localization precision is to group molecules detected multiple times  
327 in adjacent frames, and average their position<sup>36</sup>. Temporal grouping was used by the top performers  
328 (including SMolPhot<sup>37</sup>, MIATool<sup>38</sup> and SMAP-2018<sup>20</sup>), and is visibly apparent as a more punctate super-  
329 resolved image (Fig 4).

#### 330 *Choice of PSF model*

331 Most software used a variant of Gaussian PSF model. A few participants designed more accurate PSF  
332 models (Table S1). Either diffraction theory was used (MIATool<sup>38</sup>, LEAP<sup>35</sup>) or spline fitting of an  
333 analytical function to the experimental PSF was adopted (CSpline<sup>39</sup>, SMAP-2018<sup>20</sup>). Although simple  
334 Gaussian model PSFs were sufficient to obtain best-in-class performance for the 2D and astigmatic  
335 modalities (ADCG<sup>33</sup>, PeakFit, SMolPhot), top results for the more optically complex biplane and double  
336 helix modalities were exclusively PSF-modelling algorithms (SMAP-2018, CSpline, MIATool, LEAP).

#### 337 *Multi-algorithm packages*

338 Several software packages take a Swiss army knife approach of integrating multiple optional  
339 localization algorithms into one program, to be flexible enough to suit various experimental  
340 conditions<sup>20,29</sup>. SMAP and ThunderSTORM achieved strong across-the-board performance supporting  
341 this rationale.

## 342 **DISCUSSION**

343 We performed the first broad evaluation of software for 3D single molecule localization microscopy,  
344 to assess the state of the field and to allow non-specialists to determine the optimal software for their  
345 experiments.

346 In order to provide a realistic assessment of 3D software performance we tested software on  
347 simulations incorporating experimentally acquired microscope point spread functions. Our  
348 experimental-PSF-derived simulation approach is readily adaptable to novel engineered 3D SMLM  
349 PSFs<sup>40</sup> or to the PSF of individual microscopes. For instance, it would be possible to combine our  
350 derived-PSF approach with the SMLM sample simulation tool SuReSim<sup>41</sup> in order to generate ultra-  
351 realistic synthetic data, which could then be personalized to each experimentalists sample and  
352 microscope, to easily determine the blocker factors to maximal resolution, for a given experiment.

353 The strongest conclusion we draw from the 3D localization microscopy challenge is that choice of  
354 localization software greatly affects the quality of final super-resolution data, even at “easy” high SNR,  
355 low spot density conditions. Biplane performance was particularly dependent on software choice, with  
356 only one software (SMAP-2018<sup>20</sup>) achieving near-Cramer-Rao lower bound performance. Double helix  
357 SMLM showed much less sensitivity to choice of software than biplane, and showed poorer  
358 performance overall, with astigmatic SMLM intermediate between the two. The best software in each  
359 modality performed close to the Cramer-Rao lower bounds over a wide focal range and successfully  
360 detected most molecules, even at low signal to noise. Average software in all three modalities was  
361 significantly worse, with the obtained axial resolution being particularly sensitive to software choice.

362 The second major conclusion of the 3D challenge is that localization software that explicitly includes  
363 the experimental PSF in the fitting model gives a significant performance increase for 3D SMLM. For  
364 the more optically complex biplane and double helix modalities in particular, the best results were  
365 exclusively from software using PSF modelling approaches (SMAP, CSpline, MIATool). This result also



366 highlights the need for experimental PSF modelling not only in SMLM software, but also emphasizes  
367 the high degree of experimental realism required of SMLM simulations. The clear performance  
368 advantage of experimental PSF modelling software in the 3D software challenge would have been  
369 entirely unobservable had it been run with a simple analytic PSF.

370 Of the different algorithm classes, well-tuned single-emitter and multi-emitter fitting algorithms (each  
371 capable of dealing well with occasional molecule overlap) gave good results for low density 3D SMLM.  
372 We also found that several software packages for astigmatic or biplane imaging gave adequate  
373 performance for the challenging case of high molecule densities, as long as the image SNR was high.  
374 Current software packages gave poor performance when molecule density was high and image SNR  
375 was low. These results suggest that, at least with current algorithms, high density 3D SMLM  
376 performance is mediocre at high SNR, and poor at low SNR. Surprisingly, multi-emitter fitting did not  
377 show significant improvement over well-tuned single emitter-fitting for 3D high-density; this may  
378 indicate that significant potential for improvement remains in this category.

379 The second 2D localization microscopy challenge provided the opportunity to reassess the state of the  
380 field. The performance of best-in-class 2D software over a range of conditions, at both high and low  
381 spot density, is excellent. The performance of the best-in-class software at high spot density (ADCG<sup>33</sup>)  
382 was only moderately decreased compared with the low spot density results, with nearly identical  
383 molecule detection performance, and a 30 % increase in localization error. Interestingly, the top three  
384 performers in the 2D high density condition were all compressed sensing algorithms (ADCG<sup>33</sup>,  
385 FALCON<sup>32</sup>, SMfit). In low density 2D conditions, the best single-emitter, multi-emitter and compressed  
386 sensing algorithms all gave comparable, excellent, performance. We speculate that performance in  
387 this category may now be near optimal levels.

388 Rapid improvements in sCMOS camera technology mean that these cameras are rapidly becoming a  
389 major platform for single molecule localization microscopy<sup>6</sup>. Therefore, a key future goal for SMLM  
390 software assessment should be to include sCMOS-specific localization microscopy software.  
391 Furthermore, there remain two important classes of super-resolution microscopy for which software  
392 performance is crucial, but no broad software assessment has yet been performed: fluorescence-  
393 fluctuation-based super-resolution microscopies (*e.g.*, 3B<sup>42</sup>, SOFI<sup>43</sup>, SSRF<sup>44</sup>) and structured illumination  
394 microscopy<sup>45</sup>.

395 The results of this competition clearly demonstrate the formidable algorithmic performance of the  
396 best 2D and 3D localization microscopy software. However, a key outstanding challenge that often  
397 hinders adoption of new algorithms is that only a small subset of algorithms are packaged in, or  
398 compatible with fast, well-maintained, user-friendly software packages, which include all stages of the  
399 SMLM data analysis pipeline – analysis, visualization and quantification. One solution would be for the  
400 SMLM software community to collectively adopt both a standard data format and a single software  
401 platform for future software development, such as FIJI/ ImageJ<sup>46</sup>. Any new algorithm released in this  
402 environment could be immediately and widely adopted by users, and easily integrated into existing  
403 packages for SMLM analysis, visualization and quantification.

404 Both the 3D and 2D localization challenges remain open and continuously updated on the competition  
405 website. This continuously evolving analysis of state of the art super-resolution software performance  
406 provides a valuable resource to super-resolution microscopists, helping to ensure they use software  
407 that gets the best out of hard-won data. It also provides SMLM software developers with a robust  
408 means of benchmarking new algorithms against current state of the art.

## 409 **ACKNOWLEDGEMENTS**

410 *Authors acknowledge the following funding sources: a Newcastle University Research Fellowship, and*  
411 *a Sir Henry Dale Fellowship jointly funded by the Wellcome Trust and the Royal Society (Grant Number*  
412 *206670/Z/17/Z) to SH; an European Research Council (ERC) under the European Union's Horizon 2020*

413 *research and innovation programme, Grant Agreement no. 692726 “GlobalBioIm: Global integrative*  
414 *framework for computational bio-imaging.” to DS, TAP, MU; UK BBSRC grants BB/M022374/1,*  
415 *BB/P027431/1, BB/R000697/1 grant and MRC grants MC-UU-12018/2, MR/K015826/1 to RH; and*  
416 *European Research Council (ERC) grant CoG-724489, CellStructure to JR. We thank all the localization*  
417 *microscopy challenge participants for their contribution: H. Babcock, F. Hauser, S. Watanabe, N. Boyd,*  
418 *J. Min, K. Jin, H. Rouault, E. Soubies, A. von Diezmann, C. Bayas, W.E. Moerner, J. Min, J.C. Ye, T.*  
419 *Vomhof, J. Reichel, H. Pan, Z. Huang, Y. Wang, R Belmurugan, A.V. Abraham, R. J. Ober, A. Herbert, K.*  
420 *Martens, J. Hohlbein, L. Li, R. Henriques, G. Tamas, J. Sinko, M. Kirchgessner, F. Gruell, Y. Li, J. Ries, H.*  
421 *Iloma, M. Pars, A. Loot, Y. Jung, N. Fakhri, A. Archetti, M. Ovesny, G. Hagen, P Krizek, J. Huang, A.*  
422 *Kechkar, J.B. Sibarita. We thank the SMLMS 2016 organizers (Professor S. Manley and Professor A.*  
423 *Radenovic, EPFL) for hosting a localization microscopy challenge special session. We also thank Double*  
424 *Helix LLC and Molecular Devices LLC for sponsoring the SMLMS 2016 special session. The sponsors had*  
425 *no input or influence on the research.*

## 426 **AUTHOR CONTRIBUTIONS**

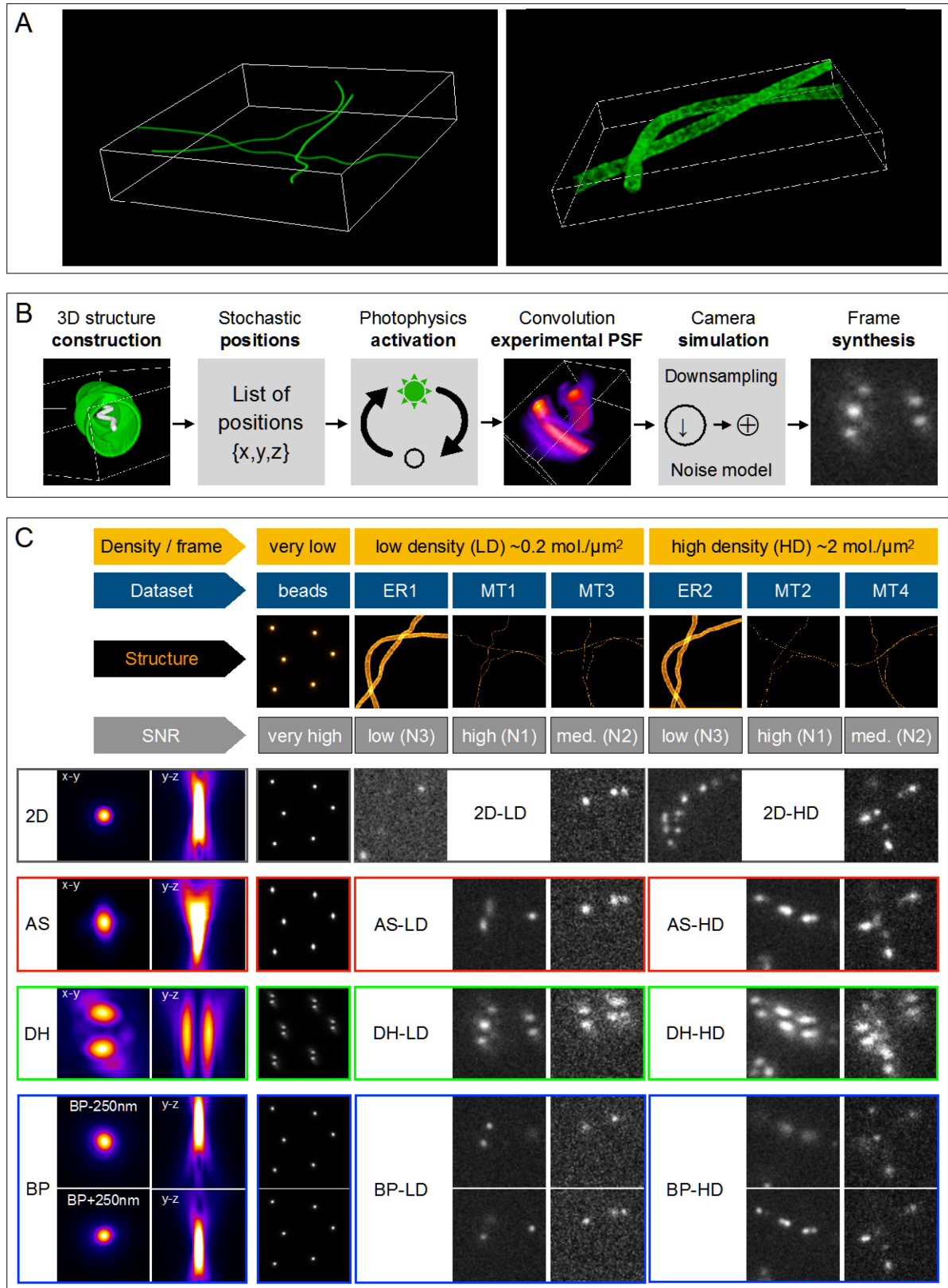
427 DS and SH conceived and coordinated the study. DS, SH, TAP, AAr, HB, SC, AW, GH, RH, TL, TP, JBS  
428 designed the study. SH, AAg, RH, JBS collected experimental PSFs. DS, TAP SH, TL wrote simulation  
429 code. BR shared unpublished software. DS generated simulated datasets. AH, JR, RV provided  
430 feedback and quality control on simulations and analysis methods. TAP carried out the assessment of  
431 software performance. TAP, DS, SH analysed and interpreted the results. DS, HB, RO, BR, GH, JBS, JR,  
432 RH, MU, SH directed research. SH, DS, TAP wrote the manuscript with feedback from all authors.

## 433 **REFERENCES**

- 434 1. Betzig, E. *et al.* Imaging Intracellular Fluorescent Proteins at Nanometer Resolution. *Science* **313**,
- 435 1642–1645 (2006).
- 436 2. Hess, S. T., Girirajan, T. P. K. & Mason, M. D. Ultra-High Resolution Imaging by Fluorescence
- 437 Photoactivation Localization Microscopy. *Biophys. J.* **91**, 4258–4272 (2006).
- 438 3. Rust, M. J., Bates, M. & Zhuang, X. Sub-diffraction-limit imaging by stochastic optical reconstruction
- 439 microscopy (STORM). *Nat Methods* **3**, 793–795 (2006).
- 440 4. Holden, S. J., Uphoff, S. & Kapanidis, A. N. DAOSTORM: an algorithm for high- density super-
- 441 resolution microscopy. *Nat Meth* **8**, 279–280 (2011).
- 442 5. Huang, F., Schwartz, S. L., Byars, J. M. & Lidke, K. A. Simultaneous multiple-emitter fitting for single
- 443 molecule super-resolution imaging. *Biomed. Opt. Express* **2**, 1377–1393 (2011).
- 444 6. Huang, F. *et al.* Video-rate nanoscopy using sCMOS camera-specific single-molecule localization
- 445 algorithms. *Nat. Methods* **10**, 653–658 (2013).
- 446 7. Sage, D. *et al.* Quantitative evaluation of software packages for single-molecule localization
- 447 microscopy. *Nat. Methods* **12**, 717–724 (2015).
- 448 8. Huang, B., Jones, S. A., Brandenburg, B. & Zhuang, X. Whole-cell 3D STORM reveals interactions
- 449 between cellular structures with nanometer-scale resolution. *Nat Meth* **5**, 1047–1052 (2008).
- 450 9. Shtengel, G. *et al.* Interferometric fluorescent super-resolution microscopy resolves 3D cellular
- 451 ultrastructure. *Proc. Natl. Acad. Sci.* **106**, 3125–3130 (2009).
- 452 10. Huang, B., Wang, W., Bates, M. & Zhuang, X. Three-Dimensional Super-Resolution Imaging by
- 453 Stochastic Optical Reconstruction Microscopy. *Science* **319**, 810–813 (2008).
- 454 11. Pavani, S. R. P. *et al.* Three-dimensional, single-molecule fluorescence imaging beyond the
- 455 diffraction limit by using a double-helix point spread function. *Proc. Natl. Acad. Sci.* **106**, 2995–2999
- 456 (2009).
- 457 12. Juetten, M. F. *et al.* Three-dimensional sub-100 nm resolution fluorescence microscopy of thick
- 458 samples. *Nat. Methods* **5**, 527–529 (2008).
- 459 13. Shtengel, G. *et al.* Interferometric fluorescent super-resolution microscopy resolves 3D
- 460 cellular ultrastructure. *Proc. Natl. Acad. Sci.* **106**, 3125–3130 (2009).

- 461 14. Babcock, H., Sigal, Y. M. & Zhuang, X. A high-density 3D localization algorithm for stochastic  
462 optical reconstruction microscopy. *Opt. Nanoscopy* **1**, 1–10 (2012).
- 463 15. Ovesný, M., Křížek, P., Švindrych, Z. & Hagen, G. M. High density 3D localization microscopy  
464 using sparse support recovery. *Opt. Express* **22**, 31263–31276 (2014).
- 465 16. Min, J. *et al.* 3D high-density localization microscopy using hybrid astigmatic/ biplane imaging  
466 and sparse image reconstruction. *Biomed. Opt. Express* **5**, 3935–3948 (2014).
- 467 17. Zhang, S., Chen, D. & Niu, H. 3D localization of high particle density images using sparse  
468 recovery. *Appl. Opt.* **54**, 7859–7864 (2015).
- 469 18. Annibale, P., Vanni, S., Scarselli, M., Rothlisberger, U. & Radenovic, A. Quantitative Photo  
470 Activated Localization Microscopy: Unraveling the Effects of Photoblinking. *PLOS ONE* **6**, e22678  
471 (2011).
- 472 19. Collaboration through competition. *Nat. Methods* **11**, 695 (2014).
- 473 20. Li, Y. *et al.* Real-time 3D single-molecule localization using experimental point spread  
474 functions. *Nat. Methods* (2018). doi:10.1038/nmeth.4661
- 475 21. Carlini, L. & Manley, S. Live Intracellular Super-Resolution Imaging Using Site-Specific Stains.  
476 *ACS Chem. Biol.* **8**, 2643–2648 (2013).
- 477 22. Shim, S.-H. *et al.* Super-resolution fluorescence imaging of organelles in live cells with  
478 photoswitchable membrane probes. *Proc. Natl. Acad. Sci.* **109**, 13978–13983 (2012).
- 479 23. Single-Molecule Localization Microscopy • Software Benchmarking. Available at:  
480 <http://bigwww.epfl.ch/smlm/challenge2016/index.html?p=results>. (Accessed: 15th June 2018)
- 481 24. Baddeley, D. & Bewersdorf, J. Biological Insight from Super-Resolution Microscopy: What We  
482 Can Learn from Localization-Based Images. *Annu. Rev. Biochem.* **87**, 965–989 (2018).
- 483 25. Fox-Roberts, P. *et al.* Local dimensionality determines imaging speed in localization  
484 microscopy. *Nat. Commun.* **8**, 13558 (2017).
- 485 26. Henriques, R. *et al.* QuickPALM: 3D real-time photoactivation nanoscopy image processing in  
486 ImageJ. *Nat Meth* **7**, 339–340 (2010).
- 487 27. Takeshima, T., Takahashi, T., Yamashita, J., Okada, Y. & Watanabe, S. A multi-emitter fitting  
488 algorithm for potential live cell super-resolution imaging over a wide range of molecular densities.  
489 *J. Microsc.* **0**,
- 490 28. Kechkar, A., Nair, D., Heilemann, M., Choquet, D. & Sibarita, J.-B. Real-Time Analysis and  
491 Visualization for Single-Molecule Based Super-Resolution Microscopy. *PLOS ONE* **8**, e62918 (2013).
- 492 29. Ovesný, M., Křížek, P., Borkovec, J., Švindrych, Z. & Hagen, G. M. ThunderSTORM: a  
493 comprehensive ImageJ plug-in for PALM and STORM data analysis and super-resolution imaging.  
494 *Bioinformatics* **30**, 2389–2390 (2014).
- 495 30. Soubies, E., Blanc-Féraud, L. & Aubert, G. A Continuous Exact  $\ell_0$  Penalty (CELO) for Least  
496 Squares Regularized Problem. *SIAM J. Imaging Sci.* **8**, 1607–1639 (2015).
- 497 31. Babcock, H. P., Moffitt, J. R., Cao, Y. & Zhuang, X. Fast compressed sensing analysis for super-  
498 resolution imaging using L1-homotopy. *Opt. Express* **21**, 28583–28596 (2013).
- 499 32. Min, J. *et al.* FALCON: fast and unbiased reconstruction of high-density super-resolution  
500 microscopy data. *Sci. Rep.* **4**, 4577 (2014).
- 501 33. Boyd, N., Schiebinger, G. & Recht, B. The Alternating Descent Conditional Gradient Method  
502 for Sparse Inverse Problems. *SIAM J. Optim.* **27**, 616–639 (2017).
- 503 34. Huang, J., Sun, M. & Chi, Y. Super-resolution image reconstruction for high-density 3D single-  
504 molecule microscopy. in *2016 IEEE 13th International Symposium on Biomedical Imaging (ISBI)*  
505 241–244 (2016). doi:10.1109/ISBI.2016.7493254
- 506 35. Pan, H., Simeoni, M., Hurley, P., Blu, T. & Vetterli, M. LEAP: Looking beyond pixels with  
507 continuous-space EstimAtion of Point sources. *Astron. Astrophys.* **608**, A136 (2017).
- 508 36. Durisic, N., Laparra-Cuervo, L., Sandoval-Álvarez, Á., Borbely, J. S. & Lakadamyali, M. Single-  
509 molecule evaluation of fluorescent protein photoactivation efficiency using an in vivo  
510 nanotemplate. *Nat. Methods* **11**, 156–162 (2014).

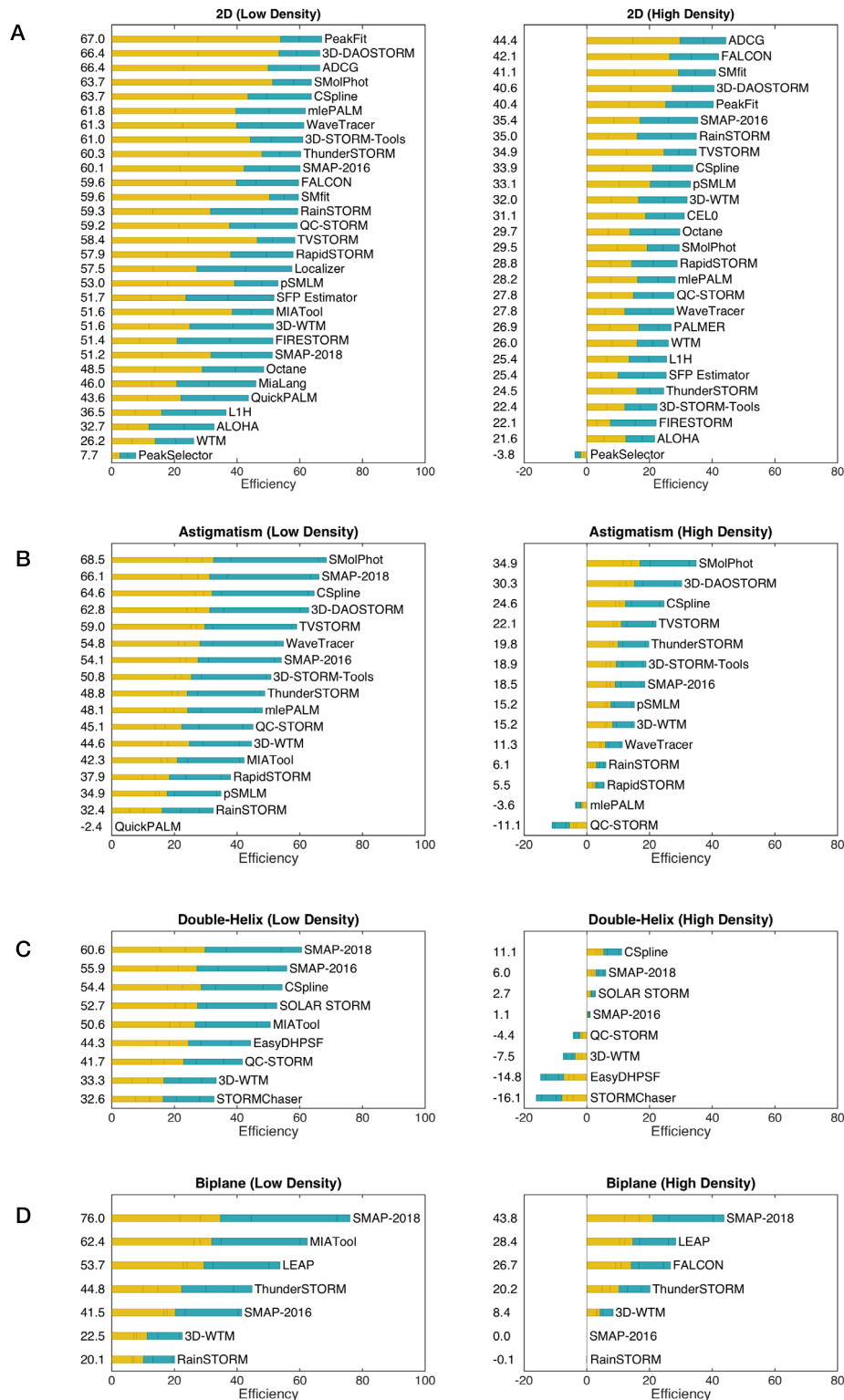
- 511 37. Martti Pars, A. L. SMolPhot Software. Available at: <https://bitbucket.org/ardiloot/smolphot->  
512 software/wiki/Home.
- 513 38. Chao, J., Ward, E. S. & Ober, R. J. A software framework for the analysis of complex microscopy  
514 image data. *IEEE Trans. Inf. Technol. Biomed. Publ. IEEE Eng. Med. Biol. Soc.* **14**, 1075–1087 (2010).
- 515 39. Babcock, H. P. & Zhuang, X. Analyzing Single Molecule Localization Microscopy Data Using  
516 Cubic Splines. *Sci. Rep.* **7**, 552 (2017).
- 517 40. Shechtman, Y., Weiss, L. E., Backer, A. S., Sahl, S. J. & Moerner, W. E. Precise Three-  
518 Dimensional Scan-Free Multiple-Particle Tracking over Large Axial Ranges with Tetrapod Point  
519 Spread Functions. *Nano Lett.* **15**, 4194–4199 (2015).
- 520 41. Venkataramani, V., Herrmannsdörfer, F., Heilemann, M. & Kuner, T. SuReSim: simulating  
521 localization microscopy experiments from ground truth models. *Nat. Methods* **13**, 319–321 (2016).
- 522 42. Cox, S. *et al.* Bayesian localization microscopy reveals nanoscale podosome dynamics. *Nat.*  
523 *Methods* **9**, 195–200 (2012).
- 524 43. Dertinger, T., Colyer, R., Iyer, G., Weiss, S. & Enderlein, J. Fast, background-free, 3D super-  
525 resolution optical fluctuation imaging (SOFI). *Proc. Natl. Acad. Sci.* **106**, 22287–22292 (2009).
- 526 44. Gustafsson, N. *et al.* Fast live-cell conventional fluorophore nanoscopy with ImageJ through  
527 super-resolution radial fluctuations. *Nat. Commun.* **7**, (2016).
- 528 45. Gustafsson, M. G. L. Surpassing the lateral resolution limit by a factor of two using structured  
529 illumination microscopy. SHORT COMMUNICATION. *J. Microsc.* **198**, 82–87 (2000).
- 530 46. Schindelin, J. *et al.* Fiji: an open-source platform for biological-image analysis. *Nat. Methods*  
531 **9**, 676–682 (2012).
- 532 47. Hanser B. M., Gustafsson M. G. L., Agard D. A. & Sedat J. W. Phase-retrieved pupil functions in  
533 wide-field fluorescence microscopy. *J. Microsc.* **216**, 32–48 (2004).
- 534 48. Hirsch, M., Wareham, R. J., Martin-Fernandez, M. L., Hobson, M. P. & Rolfe, D. J. A Stochastic  
535 Model for Electron Multiplication Charge-Coupled Devices – From Theory to Practice. *PLOS ONE* **8**,  
536 e53671 (2013).
- 537 49. Basden, A. G., Haniff, C. A. & Mackay, C. D. Photon counting strategies with low-light-level  
538 CCDs. *Mon. Not. R. Astron. Soc.* **345**, 985–991 (2003).
- 539 50. Carlini, L., Holden, S. J., Douglass, K. M. & Manley, S. Correction of a Depth-Dependent Lateral  
540 Distortion in 3D Super-Resolution Imaging. *PLoS ONE* **10**, e0142949 (2015).
- 541



544 **Figure 1: Summary of SMLM challenge simulations. A.** 3D rendering of microtubules and endoplasmic  
545 reticulum samples in a  $6.4 \mu\text{m} \times 6.4 \mu\text{m} \times 1.5 \mu\text{m}$  volume. **B. Key simulation steps.** The structure is  
546 constructed from 3D tubes continuously defined by three B-spline functions in the volume of interest.  
547 Membranes of the tubes are densely populated with possible positions. Fluorophores follow a 4-state  
548 photophysics model. Activations of a given frame are convolved with the experimental PSF and shot  
549 & camera noise is added. **C.** Summary of all 16 challenge datasets, calibration data and experimental  
550 PSFs. Each dataset is characterized by its structure (endoplasmic reticulum (ER) or microtubules (MT)),  
551 by its modality (2D, AS, DH, BP), its density (LD or HD) and by its SNR determined by the level of noise  
552 N1, N2, and N3. Left column: orthogonal projections of the experimentally-derived PSF. Eight  
553 categories were proposed for the challenge containing two datasets each, 2D-LD and 2D-HD, grey; AS-  
554 LD and AS-HD, red, DH-LD and DH-HD, green; BP-LD and BP-HD, blue.

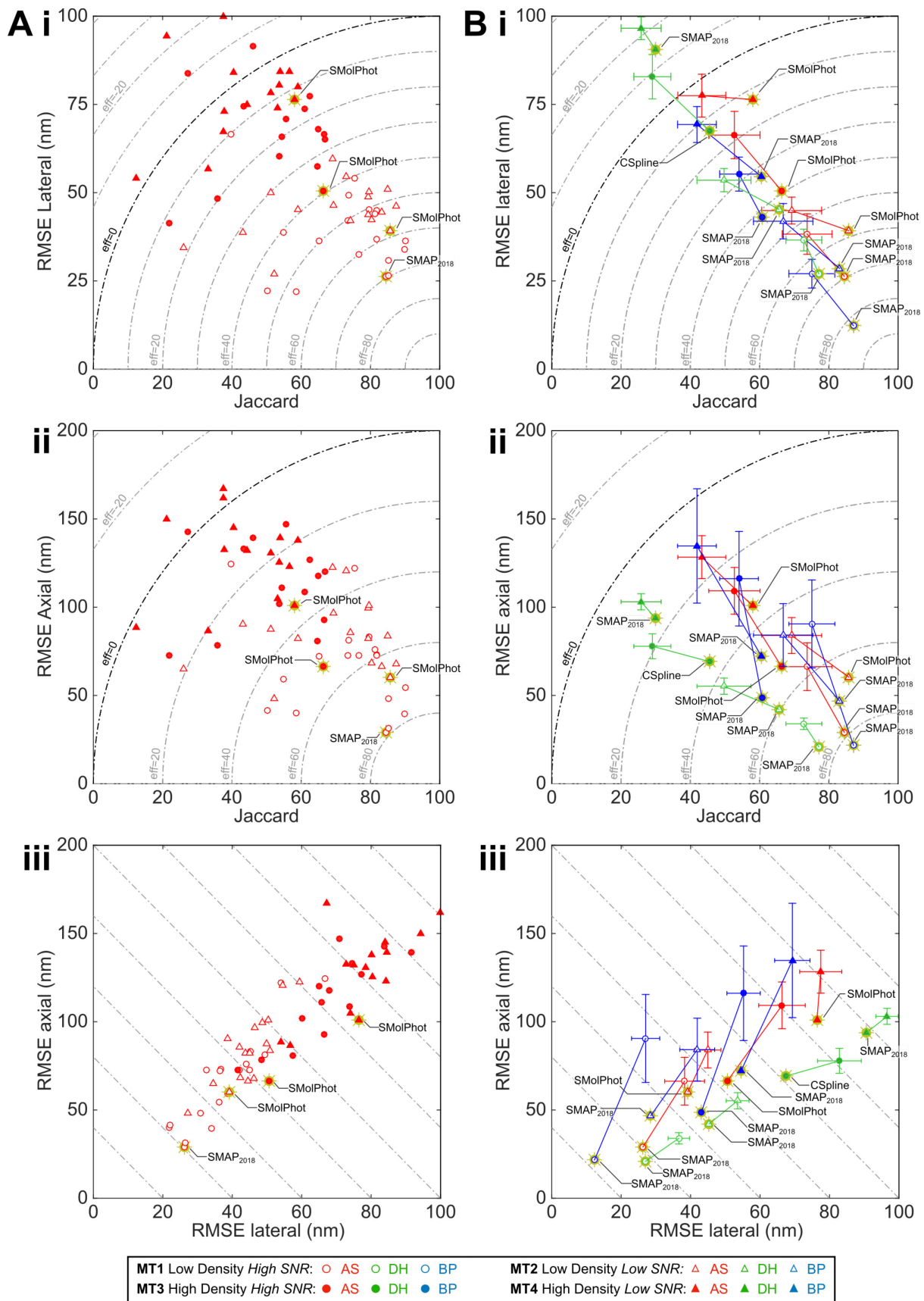
555

556



557

558 **Figure 2: Leaderboards for each competition category.** Ranking is based on the efficiency of software  
 559 based on fraction of successfully detected molecules (Jaccard index) and precision of localization  
 560 (RMSE, root mean square error, lateral & axial). The contribution of the high SNR dataset is plotted in  
 561 orange and the contribution of low SNR dataset to the efficiency is plotted in blue.



562  
563  
564  
565

**Figure 3: Comparison of 3D software performance.** Gold stars indicate top performers for each dataset. Dashed lines in top, middle panels indicate overall efficiency (higher is better). **A.** Performance of all astigmatic SMLM software (for other modality results see Supporting Material). **B.**



566 Average (colored marker with error bars) and best-in-class (colored marker with gold star) software  
567 performance for all competition modalities. *AS*, *astigmatism*; *DH*, *double helix*; *BP*, *biplane*.  
568



570 **Figure 4:** *Super-resolved images of 3D competition datasets for best-in-class (top) and representative*  
571 *average (bottom) software in each modality, for high SNR datasets. Box indicates zoomed region (left)*  
572 *or region of line profile (middle). Red, ground truth; green, software results. GT, ground truth; AS,*  
573 *astigmatism; DH, double helix; BP, biplane. Panel label key: Software\_name Ranking° (Efficiency).*

574

## 575 **METHODS**

### 576 **1. CHALLENGE ORGANIZATION**

577 We first ran the 3D SMLM software challenge as a time limited competition, with a results session  
578 hosted as a special session of the 6<sup>th</sup> Annual Single Molecule Localization Microscopy Symposium in  
579 August 2016. The competition has now been converted to a permanent software challenge accepting  
580 new submissions. Special mention to the software SMAP and 3D-WTM that participate to our eight  
581 categories (*density x modality*). The current list of participants is at:

582 <http://bigwww.epfl.ch/smlm/challenge2016/index.html?p=participants>

583 All datasets, methods, participations, and results of the challenge 2016 made available at  
584 <http://bigwww.epfl.ch/smlm/challenge2016/>. Software for simulation and analysis is hosted on the  
585 competition GitHub repository: <https://github.com/SMLM-Challenge/Challenge2016/>

### 586 **2. LOCALIZATION MICROSCOPY SIMULATIONS**

#### 587 **21. Structure**

588 The synthetic datasets were designed to be similar to images derived from cellular structures in real  
589 experimental conditions. We defined mathematical models for cellular structures that imitate  
590 cytoskeletal filaments such as microtubules and larger tubular structures such as the endoplasmic  
591 reticulum or mitochondria (Fig S1A). These structures have a tubular shape in the 3D space. Pseudo-  
592 microtubules are defined with their central axis elongating in a 3D space having an average outer  
593 diameter of 25 nm with an inner, hollow tube of 15 nm diameter. Pseudo-endoplasmic reticulum is  
594 defined as having a diameter of approximately 150nm.

595 The underlying sample structure is formalized in a continuous space which allows rendering of digital  
596 images at any scale, from very high resolution (up to 1 nm/pixel) to low resolution (camera resolution:  
597 100 nm/ pixel). The continuous-domain 3D curve is represented by means of a polynomial spline. The  
598 sample is imaged in a  $6.4 \times 6.4 \mu\text{m}^2$  field of view, and the center lines of the microtubules have limited  
599 variation along the *z* (vertical) axis, *i.e.*, less than 1.5  $\mu\text{m}$ . The fluorescent markers are uniform  
600 randomly distributed over the structure according to the required density. The photon emission rate  
601 of each fluorophore is controlled by a photo-activation model (see below).

602 The exact locations of all fluorophores are stored at high precision floating-point numbers expressed  
603 in nanometers. This ground-truth file is useful for conducting objective evaluations without human  
604 bias.

#### 605 **2.2. Photophysics activation model**

606 Given a list of source locations from the structure simulator, fluorophore blinking was simulated by a  
607 4-states Markov chain model. The states are ON, OFF, BLEACH, DARK and the transition are Poisson  
608 distributed (Fig S1C), except for the OFF to ON transitions which follow a uniform random distribution,  
609 to reflect that in typical experimental conditions, constant imaging density is maintained by tuning the  
610 photoactivation rate during the experiment. All switching is calculated at sub-frame resolution and  
611 then total fluorophore on-time was integrated over each frame.

612 Due to two decay paths, the actual mean lifetime of the state ON is

$$613 \overline{T}_{\text{LIFETIME}} = \frac{1}{\frac{1}{T_{\text{ON}}} + \frac{1}{T_{\text{BLEACH}}}}$$

614 Switching rates were chosen to approximate photoactivatable fluorescent proteins  $T_{\text{ON}} = 3$  frame,  $T_{\text{DARK}}$   
615  $= 2.5$  frames, and  $T_{\text{BLEACH}} = 1.5$  frames.

616 Fractional fluorophore ON-times per frame (between 0 and 1) were then multiplied by the mean flux  
617 of photon emission. The flux of photons expressed in photons/seconds was given by the relation

$$618 \quad \mathbf{F} = \frac{\phi \cdot P \cdot \sigma}{e}$$

619  $\Phi$  is the quantum yield of the dye, P is power of the laser in W/cm<sup>2</sup>,  $e = h c / \lambda$  is the energy of one  
620 photon,  $\sigma = 1000 \ln(10) \epsilon / N_A$  is the absorption cross section in cm<sup>2</sup> and  $\epsilon$  is the molar extinction  
621 coefficient (EC) or absorptivity in cm<sup>2</sup>/mol which is a characteristic of a given fluorophore. The laser  
622 power was Gaussian distributed over the field of view. At the end of this process a list of XY positions,  
623 on-frames and (noise-free) intensities for all activated fluorophores was obtained.

### 624 **2.3. Experimental Point-Spread Function**

625 Model PSFs, stored as high resolution look up tables, were derived from experimentally measured  
626 PSFs. Although the algorithmic approach is distinct, this concept of accurately modelling the  
627 experimental PSF based on calibration data bears relation to the PSF phase retrieval approach  
628 previously employed by Hanser and coworkers<sup>47</sup>.

629 Images of fluorescent beads were recorded for each modality (Table S4). Signal to noise ratio of  
630 recorded PSFs was maximized in all cases by maximizing exposure time and averaging over several  
631 frames to increase dynamic range.

632 To acquire experimental PSFs, we took 100 nm Tetraspek beads (Invitrogen) adsorbed to #1.5 (170  $\mu$ m  
633 thick) coverglass, imaged in water. The excitation wavelength was between 640 nm and 647 nm, and  
634 a Cy5 emission filter was used. Exact data acquisition parameters for each modality are listed in Table  
635 S4.

### 636 **2.4 Simulation PSF construction**

637 For each modality, 3-6 beads were selected within a small (< 32  $\mu$ m) region, to minimize PSF variation  
638 due to spherical aberration. Images for each selected bead were interpolated in XY to a pixel size of  
639 10 nm. Beads were then coaligned by cross-correlation on the in-focus frame. Coaligned beads were  
640 averaged in XY to minimize pixel quantization artefacts and to increase SNR. Where necessary, Z-stacks  
641 were interpolated to a Z-step size of 10 nm. A central Z-range of 1.5  $\mu$ m was selected that represents  
642 151 optical planes with a Z-step of 10 nm. The Z-range covers -750 nm to +750 nm. The plane of best  
643 focus was chosen as the simulation 0 nm plane. Each model PSF was normalized such that the total  
644 intensity of the PSF in the in-focus frame within a diameter of 3 FWHM from the PSF center was equal  
645 to 1.

646 For the DH PSF, the transmission of the combined phase mask system was measured as 96 %, which  
647 was approximated as 100 % brightness relative to the 2D and astigmatic PSFs.

648 In biplane super-resolution microscopy, emitted fluorescence is split into two simultaneously imaged  
649 channels, with a small (500-1000 nm) defocus introduced between the two channels<sup>12</sup>. As the small  
650 defocus should introduce minimal additional aberration into an optical system, we semi-synthetically  
651 constructed a realistic biplane PSF from the experimental 2D PSF. The two defocused PSFs were  
652 constructed by duplicating the 2D PSF and offsetting it by -250 nm and 250 nm for each Z-plane.

653 This yielded five high SNR model PSFs with an isotropic voxel size of 10x10x10 nm<sup>3</sup>. These normalized  
654 PSFs are provided on the competition website: <http://bigwww.epfl.ch/smlm/challenge2016/psf>

655 The ground truth XY=0 was defined as the image centre of mass of the in-focus frame of the model  
656 PSF, and Z=0 was defined as the in-focus frame. Accounts for shifts in the fitted XY centre of the model  
657 PSF by localization software due to systematic offsets and Z-dependent variation of the model PSF  
658 centre of mass are dealt with below (wobble correction).

## 659 2.4. Noise model

660 A constant mean autofluorescent background was added to the noise-free simulated images, and  
661 these images were then fed through the noise model representing Poisson distributed fluorescence  
662 emission recorded on a high quantum efficiency back-illuminated EMCCD<sup>48,49</sup>.

663 The proposed noise model assumed as main contributions to the stochastic noise:

- 664 •  $\sigma_S$ , the shot noise produced by the fluorescence background and signal and the spurious  
665 charge. Shot noise can be derived from the second moment of the Poisson distribution
- 666 •  $\sigma_R$ , the read noise of EMCCD camera, which is described by second moment of the Gaussian  
667 distribution
- 668 •  $\sigma_{EM}$ , the electron multiplication noise introduced by the gain process, which is described by  
669 the second moment of the Gamma distribution<sup>49</sup>.

670

671 We assumed as camera parameters the ones specified for the Photometrics Evolve Delta 512 EMCCD  
672 camera:

- 673 • QE = 0.9, Evolve quantum efficiency at 700 nm absorption wavelength.
- 674 •  $\sigma_R = 74.4$  electrons, manufacturer measured root mean square noise for Evolve 512 camera
- 675 •  $c = 0.002$  electrons, manufacturer quoted spurious charge (clock induced charge only, dark  
676 counts negligible)
- 677 •  $EM_{gain} = 300$
- 678 •  $e_{adu} = 45$  electron per analog to digital unit (ADU), analog to digital conversion factor
- 679 •  $G = 0.9 \cdot 300 / 45 = 6$ , total system gain
- 680 • BL = 100 ADU

681 The final simulated photon electrons will thus be given by:

$$682 \quad n_{ie} = \mathcal{P}(QE \cdot n_{photIn} + c)$$

$$683 \quad n_{oe} = \Gamma(n_{ie}, EM_{gain}) + \mathcal{G}(0, \sigma_R)$$

684 which leads to the final pixel counts:

$$685 \quad ADU_{out} = \min\left(\frac{n_{oe} - n_{oe} \bmod e_{ADU}}{e_{per\ adu}} + BL, 65535\right)$$

## 686 2.5. Depth-dependent lateral distortion: Wobble

687 As the PSF models are experimentally derived, the 3D estimated localizations exhibit a depth-  
688 dependent lateral distortion, here called *wobble*. This optical distortion is due to a combination of a  
689 systematic offset (arbitrary definition of PSF center) and optical aberrations<sup>50</sup>. In order to compare  
690 estimated and true localizations, we correct this effect during the assessment (Section 3.1).

## 691 2.5 Comparison of software results between different modalities.

692 The intensities of the PSF in each imaging modality were normalized to facilitate comparison of results  
693 between different modalities. Software results between 2D, 3D AS and 3D DH modalities are expected  
694 to be directly comparable.

695 For the biplane model PSF, as the emitted fluorescence is split into two channels, the intensity in each  
696 of the two simulated biplane channels was additionally reduced by 50 %. We note that the  
697 fluorescence background was not reduced by 50 % as intended, leading to artificially high background  
698 for the biplane simulation (*i.e.*, the background in each biplane channel is the same as in the single  
699 channel of the other modalities). However, due to the low background level in the 3D simulations, the  
700 effect on image SNR and thus localization error is small (see Fig S7), less than 5nm near the plane of

701 focus. Therefore, as long as the small drop in image SNR is taken into account, approximate  
702 comparisons of the biplane data to the other modalities can still be made.

### 703 **3. SOFTWARE ASSESSMENT**

#### 704 **3.1 Protocol**

705 Each localization file submitted by the participants was manually checked for erroneous systematic  
706 errors in the definition of the dataset coordinate system, such as offsets, XY axis flips or clear scaling  
707 errors. Datasets were then programmatically standardized into a consistent output format. All  
708 modifications are publicly available. If required, the modifications consisted of columns reordering,  
709 reversing axes, XY axis swap, and shifting the lateral positions by a half camera pixel.

710 The assessment pipeline includes three main parts: localization processing, the pairing between true  
711 and estimated localization and the metrics calculations. The first one depends on the assessment  
712 settings. There are two switchable properties: photon thresholding and wobble correction. Their  
713 combinations yield four different assessment settings. Up to 64 assessment runs per software were  
714 possible (*i.e.*, 4 modalities, 4 datasets per modality). For any setting, we excluded the fluorophores  
715 within a lateral distance of 450 nm from the border. This value corresponds to the radius of the largest  
716 PSF (*i.e.*, Double Helix). The activations too close from the border are more difficult to localise and  
717 could bias the results.

718 The pairing between true and estimated localizations was performed frame by frame. The procedure  
719 matches two sets of localizations. We deployed the presorted nearest-neighbor search for its  
720 efficiency. The results are effectively similar to the computationally intensive Hungarian algorithm<sup>7</sup>.

#### 721 *Photon thresholding*

722 A photon threshold was required primarily due to the use of a realistic fluorophore blinking model.  
723 Since a fluorophore could activate/ bleach at any point in a simulated frame, this led to many frames  
724 containing very dim, undetectable localizations, eg. where a molecule had been active for one or more  
725 frames previously, and then bleached during the first 5 % of a frame. These fractional localizations  
726 should also be present but practically undetectable in an experimental dataset.

727 In order to focus the software analysis on the localizations where the molecule was active for the  
728 majority of a frame, which we decided was most consistent with experimental expectations, we  
729 implemented a photon threshold means where we kept the 75% brightest ground truth fluorophore  
730 activations. Because this was performed *after* the pairing step, observed localizations that were paired  
731 to discarded ground truth activations were also removed from the metric calculations.

#### 732 *Wobble correction*

733 The centroid of experimental point spread functions shifts laterally by as much as 50 nm, as a function  
734 of axial position<sup>10,50</sup>. This is most often ignored by localization software, and instead corrected post-  
735 hoc by reference to a calibration curve<sup>37</sup>. Since our simulated PSF is experimentally derived, it was  
736 necessary to correct for these artefactual shifts between the observed localizations and ground truth,  
737 as part of the assessment process. This correction was performed using calibration data uploaded by  
738 competitors, similar to the correction typically performed on experimental data<sup>50</sup>.

739 Three scenarios were proposed to the participants: no correction was applied during the assessment;  
740 the correction was based on a file provided by the participant itself or the correction was calculated  
741 by ourselves. The latter nevertheless requires the participant to localize a stack of beads we provided.  
742 Since the true positions of the beads are known, the difference between the estimated and true  
743 positions could be calculated and averaged. It thus yields the values for wobble correction.

744 In certain specific cases (identified on the competition website), at the request of authors, we did not  
745 apply this correction, for example because the software explicitly considered the whole 3D PSF during

746 fitting and was thus immune to this lateral shift artefact. For accurate results, application of lateral  
747 shift correction is critical for analysis of localization microscopy simulations using experimentally  
748 derived PSFs, as can be seen by comparison of typical software results with and without wobble  
749 correction (Fig S11).

### 750 **3.2 Metrics**

751 The metrics are split into two categories: localization based and image based metrics.

752 The former directly relies on the localizations positions and notably includes the Recall, the Precision,  
753 the Jaccard Index, the RMSE (axial and lateral) and the consolidated Z-range. For the calculation of  
754 average software performance (Fig 3B) outlier software with an efficiency less than  $eff=30$  were  
755 excluded from the measurement.

756 The image based metrics are computed from a rendered image and includes the Signal-to-Noise Ratio  
757 (SNR) and the Fourier Ring / Shell Correlation (FRC/FSC). To render the image, we added the  
758 contribution of each localized molecule at the corresponding pixels. A contribution takes the form of  
759 a 3D additive Gaussian with a Full-Width Half Maximum (FWHM) of 20 nm. A complete list of all  
760 computed metrics is shown in the Supplementary Note 2.

761 We also calculated localization based metric results as a function of axial position. We proceeded by  
762 considering a subset of activations lying within an interval of axial positions (*i.e.*, from the true  
763 localizations). Then, most of the metrics (*e.g.*, Recall) are locally computed. This yields a curve  
764 providing information on the depth performance of each software / modality.

765 In order to summarize software axial performance, we analyzed how the recall varied as a function of  
766 Z. A typical recall versus axial position curve (Fig S9) will drop at positions far from the focal plane,  
767 *i.e.*, where software can no longer detect spots to defocus. We first smoothed the curve using a sliding  
768 window. Then we computed the software Z-range, defined as the full width half maximal Recall of the  
769 smoothed curve (Fig S12). This quantity is visually intuitive and useful for discussion of the recall  
770 performance if considered alongside a plot of recall vs axial position. However, because FWHM recall  
771 depends on the maximal recall, ranking based on this procedure would promote a software which  
772 poorly performed everywhere (*i.e.*, flat curve), whereas a software which performed well in the focal  
773 plane but less well outside would obtain a worse FWHM recall. This observation leads us to produce  
774 a so-called consolidated Z-range, by multiplying the Z-range value by the maximal Recall, which should  
775 provide a robust metric that avoids the previous case scenario.

776 *Principal component analysis.* In order to analyse the relationship between analysis metrics we  
777 computed the covariance matrix between each metric and the principal component analysis (PCA) on  
778 the metrics (Fig S14B). Each metric was standardized before applying the covariance and the PCA. For  
779 convenience, we took the additive inverse of the metrics for which lower values are best (*i.e.*, FP, FN,  
780 RMSE, FRC, FSC).

781 Summary statistics and detailed results for each software are available on the competition website  
782 (<http://bigwww.epfl.ch/smlm/challenge2016/index.html?p=results>), which also includes a tool for  
783 side-by-side comparison of the results of multiple software packages

### 784 **3.3 Baseline Localization Software**

785 We developed a minimalist Java tool software that performs localizations of bright emitters on the 4  
786 modalities of the challenge 2016: 2D, Astigmatism, Double-Helix, and Biplane. This  
787 SMLM\_BaselineLocalization software is only designed to establish the performance baseline for the  
788 SMLM challenge. It has intentionally limited lines of code and relies only on few threshold parameters  
789 to localize particles. It has basic calibration tool that has to run on a z-stack of beads to find the linear  
790  $f(x)$  relation between the axial position Z and the shape of the bead.



- 791       • Astigmatism:  $Z = f(W_x - W_y)$ , where  $W_x$  and  $W_y$  are respectively an estimation of the size in X  
792       and Y.
- 793       • Double-Helix:  $Z = f(\theta)$ , where  $\theta$  is the angle formed the pairing of two close points.
- 794       • Biplane:  $Z = f(W_{\text{left}} - W_{\text{right}})$ , where  $W_{\text{left}}$  and  $W_{\text{right}}$  are respectively an estimation of the size of  
795       the spots in left and the right plane.
- 796       The Java code is available: <https://github.com/SMLM-Challenge/Challenge2016>